

## 分散共有メモリ型マルチプロセッサ「ASURA」の 階層性とその評価

城 和貴<sup>†</sup>, 柳原 守<sup>†</sup>, 田中高士<sup>†</sup>, David FRASER<sup>†</sup>, 新田 博之<sup>†</sup>  
齋藤 秀樹<sup>‡</sup>, 森 眞一郎<sup>‡</sup>, 富田眞治<sup>‡</sup>  
阿草清滋<sup>◇</sup>

<sup>†</sup>: (株)クボタ    <sup>‡</sup>: 京都大学 工学部,    <sup>◇</sup>: 名古屋大学 工学部,

分散共有メモリ型マルチプロセッサ「ASURA」の性能評価を確率モデルを利用して行なう。モデル化にあたっては、ASURAの特徴である階層性に着目し、1サイクル内のキャッシュ、ローカル・メモリ、グローバル・メモリに対する可能なリクエスト数をそれぞれ求め、それらを統合してシステム全体の理論性能値とする。クラスタ内の評価にはマルコフ・チェイン・モデルを、クラスタ間の評価には予想されるリクエストのキューイング数をそれぞれ組み込み、全体として統合されたときの誤差を少なくする手法を提案する。さらにそのモデルを用いてキャッシュのヒット率及びリクエストの分散度合と理論性能との関係を分析する。

### Hierarchy Impact and Performance Evaluation of ASURA: A Distributed Shared Memory Multiprocessor

Kazuki JOE<sup>†</sup>, Mamoru YANAGIHARA<sup>†</sup>, Takashi TANAKA<sup>†</sup>, David FRASER<sup>†</sup>, and  
Hiroyuki NITTA<sup>†</sup>  
Hideki SAITO<sup>‡</sup>, Shin-ichiro MORI<sup>‡</sup>, Shinji TOMITA<sup>‡</sup>  
Kiyoshi AGUSA<sup>◇</sup>

<sup>†</sup>: Office of Computer Business, KUBOTA Corporation  
ASTEM RI  
17 Chudoji, Minami-machi, Shimogyo-ku, Kyoto 600 Japan

<sup>‡</sup>: Department of Information Science  
Faculty of Engineering, Kyoto University  
Yoshida-hon-machi, Sakyo-ku, Kyoto 606-01 Japan

<sup>◇</sup>: Department of Electronic Information  
Faculty of Engineering, Kyoto University  
Yoshida-hon-machi, Sakyo-ku, Kyoto 606-01 Japan

*E-mail: {joe, yanagi, takashi, david, nitta}@kocb.astem.or.jp E-mail: {saito, moris,  
tomita}@kuis.kyoto-u.ac.jp  
E-mail: agusa@nuee.nagoya-u.ac.jp*

A new analytic performance model is presented for ASURA, a distributed shared memory multiprocessor. The system bandwidth is analyzed as the total number of requests in a system cycle. Following the fundamental hierarchy of ASURA, the requests are divided into cache, local memory and global memory types. Furthermore, we adopt the adaptive method which increases the accuracy of overall model: 1) Markov Chain Model for the analysis of intra cluster activity. 2) Queueing probability for inter cluster requests. Finally, using this model we analyze the system performance from the viewpoint of cache hit ratio and the distribution of requests.

## 1 はじめに

プロセッサ・クラスタ方式の並列計算機は将来のスーパーコンピューティングへの有力な一方式として、商用 [26]、研究用 [12][8][9] 共に研究開発が推進されている。我々は既に大規模な並列処理計算機への第一歩として、クラスタ・ベースの階層型マルチプロセッサ・システム ASURA を提案している [27]。そこで本稿では、ASURA の階層性に着目した確率モデルを検討し、性能評価を行なう。

並列計算機の性能評価のためのモデル化は単純なクロスバ結合 [5][11][14] やマルチバス [20][6] に対するものから、クラスタ・ベース [16][15][3][2][1] のものまで、さまざまな手法が提案されている。本稿で検討する確率モデルを用いたもの以外にも、マルコフ・モデル [18]、セミ・マルコフ・モデル [19]、ベトリ・ネットワーク・モデル [17]、M/G/1 モデル [21] 等さまざまな手法が提案され、アーキテクチャの定量的性能評価に利用されている。

階層的な構造を持つクラスタ・ベースの並列計算機の性能評価を検討する場合、その階層性を無視してデータのアクセス・パターンを一律と仮定すると、計算される性能は究めて悪くなることが報告されている [16]。そのため、数年前のクラスタ・ベースの並列計算機の評価モデルには、「Favorite Memory」の概念が導入された [7][15][3]。これはソフトウェアの観点から見た場合、数値計算等で自然に発生するデータの局所性、さらにその最適化 [13]、これを一般化したクラスタ間のデータ分割配置 [22][10] といった諸問題に何らかの有効な手段が施された結果と見る事が出来よう。そして、近い将来、システムの階層性に特有なデータ・アクセスの分布はパラメータとして与えることが出来るものと期待される。本稿での手法も、ASURA の階層性に適合すると思われるデータ・アクセスの分布を仮定して、性能評価を検討していく。

本稿では ASURA の概要を述べた後、モデルの仮定を示し、各階層のモデルを統合する方式でシステム全体のモデルを表し、性能評価を行なう。

## 2 ASURA の概要

ASURA とは株式会社クボタと京都大学とが共同で開発している実験システムである [27][25]。図 1 にそのアーキテクチャを示す。ASURA はメモリの階層性に重点を置いたプロセッサ・クラスタ方式の密結合型マルチプロセッサ・システムである。

プロセッシング・エレメント (PE) はオンチッ

プの 1 次キャッシュと 4MB の 2 次キャッシュを持つ R4000MC である。これらのローカル・キャッシュはスヌープ方式のキャッシュ・コヒーレンス制御によるイリノイ・プロトコル [4] を一部変更したものに従う。

8 個の PE と 256MB のローカル・メモリ、ネットワーク・インタフェース (NIF) をバスによって結合し、1 つのプロセッサ・クラスタ (PC) を構成する。このローカル・メモリはクラスタ内の PE によって共有される。NIF は 16MB のグローバル・メモリと 32MB のグローバル・キャッシュから構成される。グローバル・メモリは ASURA システムの分散共有メモリとして全ての PE によって共有される。グローバル・キャッシュはクラスタ内の PE で共有され、フルマップ・ディレクトリ方式のキャッシュ・コヒーレンス制御によるシナプス・プロトコル [4] に従う。また、グローバル・キャッシュは上位のキャッシュとの間に MLI 特性を持つため [27] 1 次及び 2 次キャッシュはグローバル・キャッシュをキャッシング出来る。

複数の PC を階層型のネットワークによって結合することにより ASURA システムが形成される。現段階では、4 つの PC をレジスタ・インサージョン・リング [23] を一部変更した方式でリング結合しクラスタ・グループを形成する。さらに 32 クラスタ・グループをクロスバ結合することにより、最大 1024 CPU を構成することが出来る。

ASURA の特徴はそれぞれ階層的なキャッシュ、メモリ、ネットワークである。キャッシュは前述したように 3 階層から成り、メモリはある PE から見てローカル・メモリ、クラスタ内のグローバル・メモリ、クラスタ・グループ内のグローバル・メモリ、それ以外のグローバル・メモリと 4 階層からなる。また、ネットワークはクラスタ内のバス、クラスタ・グループ内のリング、全体にわたるクロスバと 3 階層からなる。

## 3 準備

### 3.1 ASURA のモデル化

性能評価のモデル化のために、ASURA を次のように定義する。P 個のプロセッサ (プライベート・キャッシュ (PCH) を持つ) と 1 つの<sup>1</sup>メモリ・モジュール、ネットワーク・インタフェース (NIF) がバス結合され、プロセッサ・クラスタ (PC) を形成する。ただし、NIF はグローバル・メモリ (GM) とグローバル・キャッシュ (G

<sup>1</sup>実際にはインターリーブされているが、バスとメモリのサイクル比が同じと考えた方がモデルが簡単であるためこのような仮定を用いる

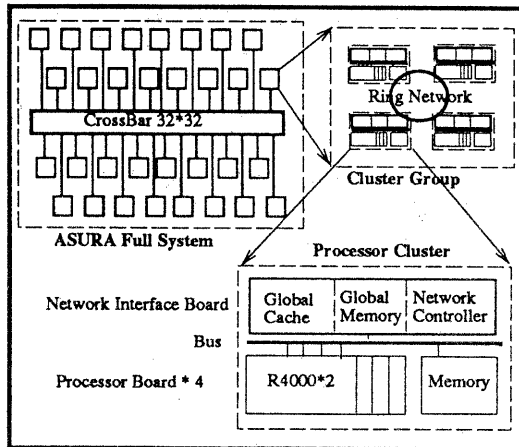


図 1: ASURAのアーキテクチャ

CH)、PC間の通信を制御するコントローラからなる。(PC間の通信はPC内のバスに影響を与えないことに注意) M個のPCはリング結合されてクラスタ・グループ(CG)を構成し、K個のCGがクロスバ結合することによりASURA全体が定義される。従って、全体では $P \times M \times K$ 個のプロセッサとPCH、それぞれ $M \times K$ 個のローカル・メモリ・モジュール、グローバル・メモリ・モジュール、GCH、 $M \times K$ 本のバス、K個のリング、1個のクロスバによって並列計算機システムを構成することになる。メモリ階層性を示すために、メモリ・ロケーションをローカル・メモリ(LM)、PC内グローバル・メモリ(PCGM)、CG内グローバル・メモリ(CGGM)、リモート・グローバル・メモリ(RGM)で区別する。

### 3.2 モデルの仮定

本稿ではASURAの階層的なメモリ及びネットワークの性能を評価するために、システム全体のバンド幅を確率的なモデル[20][14][6][7]によって求める。モデルを簡略化するために次のような仮定を導入する。

1. 全てのプロセッサからのリクエストは全て独立である
2. リクエスト先ははある確率分布関数に従う
3. アクセプトされなかったリクエストはキャンセルされる
4. ネットワークの遅延時間はリング以外は無視する
5. PCH、GCHのコヒーレント動作は考慮に入れない

仮定2はデータの局所性を意味している。キャッシュはデータの局所性を利用した機構であり、様々なアーキテクチャに利用できる。これは並列計算機におけるキャッシュでも同じであり、この性質をメモリにまで適用すると、LM、PCGM、CGGM、GMの参照確率が一樣ではなく、参照時間と参照確

率の頻度とのトレードオフによって階層的なモデルを形成出来る[24]。例えば、確率密度関数 $P_r$ を、 $Pr[X_{PCH}] > Pr[X_{GCH/LM}] > Pr[X_{PCGM}] > Pr[X_{CGGM}] > Pr[X_{GM}]$ のように設定出来る。ただし、同一階層においては参照確率は一樣であるとする。例えば、CG内のリクエストが特定のPCである確率は $1/M$ である。同様、RGMへのリクエストが特定のCGである確率は $1/K$ である。

仮定3は幾分非現実的である。なぜならば、アクセプトされなかったリクエストは次のサイクルで再発行されるからである。しかしながら、この仮定はモデルを大幅に簡単にし、実際の結果との誤差も小さいと報告されている[5]。本モデルでは、PC内のみ仮定3の条件を緩和して、マルコフ・チェイン・モデルを用いた手法[11]で、リクエスト率の誤差修正を図る。

仮定4ではリングの遅延はバスやクロスバのそれと比較して無視できないので、リングの遅延についてはモデルの対象としている。

並列計算機の定量的評価を行なう際には、評価の中心であるバンド幅の定義を明確にしなければならない。これまでに提案されている手法には、

- プロセッサとメモリのサイクルが等しく、ネットワークの遅延は無視するという仮定の上で、1サイクルに対するリクエスト数を持って全体のバンド幅とするもの[16][3]
- プロセッサとメモリのサイクル比を与え、メモリのサイクルごとにビジー状態のメモリ・モジュールの数を算出したものを持って全体のバンド幅とするもの[20]
- プロセッサのスピードを与え、そのスピードとクラスタ数でネットワークの遅延を近似し、適当なプロセッサ利用率を与えて全体の性能を算出するもの[1]

などがある。本モデルでは、ネットワークが複数個あることから、最もクリティカルな部分、すなわちPC内バスのサイクル・タイムを基準にして、プロセッサ及びその他のネットワークのサイクル・タイムの比を与え、基準のバス・サイクルでのシステム全体(プライベート・キャッシュ、ローカル・メモリ、グローバル・キャッシュ、グローバル・メモリ)のリクエスト数を持ってシステム全体のバンド幅と定義する。

### 3.3 記数法

P PC内のプロセッサ数  
M CG内のPC数

K 全体のCG数

$W_{CG}$  CG内でのリクエストの平均待ちサイクル数

L CG内の各ノードでのリクエストの最大キューイング数

BW システム全体のバンド幅

$BW_{local}$  ローカル・メモリに対するバンド幅

$BW_{global}$  グローバル・メモリに対するバンド幅

$\Psi$  プロセッサのリクエスト率

$\Psi_{LM}$  プロセッサのLMに対するリクエスト率

$\Psi_{GCH}$  プロセッサのGCHに対するリクエスト率

$\Psi_{PCGM}$  プロセッサのPCGMに対するリクエスト率

$\Psi_{CGGM}$  プロセッサのCGGMに対するリクエスト率

$\Psi_{RGM}$  プロセッサのRGMに対するリクエスト率

$\Psi_{ring}$  リングに対するリクエスト率 (リングのサイクル基準)

$SLM$  リクエスト先がLMである確率

$SPCGM$  リクエスト先がPCGMである確率

$SCGGM$  リクエスト先がCGGMである確率

$SRGM$  リクエスト先がRGMである確率

$h_{PCH}$  PCHに対するヒット率

$h_{GCH}$  GCHに対するヒット率

$\alpha$  プロセッサとバスのサイクル・タイムの比

$\beta$  バスとリング・ネットワークのサイクル・タイムの比

$P_{CG}$  リング・ネットワークにPCから入力される確率

$P_x$  CGからクロスバ・ネットワークに入力される確率

$q_{CG(i)}$  CGに*i*個のリクエストが発行される確率

$E_{CG(i)}$  CGに発行された*i*個のリクエストがアクセプトされる確率

## 4 モデル

### 4.1 PCのモデル

一般にPCHを持つP個のプロセッサが1個のメモリモジュールにバス結合されている場合を考える。プロセッサのサイクルごとのリクエスト率を $\Psi$ 、キャッシュのヒット率を $h$ 、プロセッサとバスのサイクルの比を $\alpha$ とした場合、あるプロセッサがバスの1サイクル内にリクエストを出す確率は

$$\alpha(1-h)\Psi \quad (1)$$

よって、少なくとも1つのプロセッサがバスのサイクル内にリクエストを出す確率は

$$1 - (1 - \alpha(1-h)\Psi)^P \quad (2)$$

仮定より、バスへのリクエストはそのままメモリへの参照と見なされる。また、各プロセッサはバスのサイクル毎に $\alpha h \Psi P$ 回のリクエストを自分のPCHに発行している。よって、PCのバンド幅は

$$BW_{(0)} = 1 - (1 - \alpha(1-h)\Psi)^P + \alpha h \Psi P \quad (3)$$

バスへのリクエスト数は $\alpha(1-h)\Psi P$ であるから、リクエストがアクセプトされる確率 $P_A^{(0)}$ は

$$P_A^{(0)} = \frac{BW_{(0)}}{\alpha(1-h)\Psi P} \quad (4)$$

ここで、前章仮定3の条件を緩和してリクエスト率の誤差修正を図る。ももとの仮定ではキャンセルされたリクエストは捨てられるわけだが、実際は同じリクエストを次のサイクルで再発行することになる。これを、1サイクルの待ち状態の後、即座に別のリクエストを発行すると考えると、プロセッサの状態は、リクエストを発行しているアクティブな状態(A)と、キャンセルされた待ち状態(W)の2状態に分類され、その状態推移はマルコフ・チェーンを形成する。図2はその状態推移を表す。a1,a2,a3,a4はそれぞれ、状態(A)から

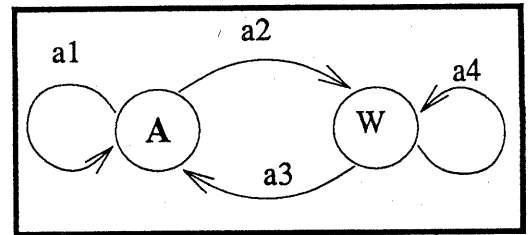


図2: 状態推移を示すマルコフ・グラフ

(A)、(A)から(W)、(W)から(A)、(W)から(W)への推移確率で、以下のように表される。(ただし、 $\psi = \alpha(1-h)\Psi$ とする。)

$$a1 = (1 - \psi) + P_A^{(0)}\psi \quad (5)$$

$$a2 = \psi(1 - P_A^{(0)}) \quad (6)$$

$$a3 = P_A^{(0)} \quad (7)$$

$$a4 = 1 - P_A^{(0)} \quad (8)$$

a1,a2,a3,a4は推移確率行列を形成するので、以下のようにして確率ベクトルを求めることができる。

$$\lim_{m \rightarrow \infty} \begin{pmatrix} a1 & a2 \\ a3 & a4 \end{pmatrix}^m = \frac{1}{P_A^{(0)} + \psi(1 - P_A^{(0)})} \begin{pmatrix} P_A^{(0)} & \psi(1 - P_A^{(0)}) \\ P_A^{(0)} & \psi(1 - P_A^{(0)}) \end{pmatrix} \quad (9)$$

よって、状態(A),(W)の確率 $Q_A, Q_W$ は

$$Q_A = \frac{P_A^{(0)}}{\psi + (1 - P_A^{(0)})} \quad (10)$$

$$Q_W = 1 - Q_A \quad (11)$$

となる。状態(A)の時のリクエスト率は $\Psi Q_A$ 、状態(W)の時のリクエスト率は $Q_W$ だから、修正されたリクエスト率 $\Psi^{(1)}$ は

$$\Psi^{(1)} = \Psi^{(0)} Q_A + Q_W \quad (12)$$

これを $P_A^{(0)}$ を使って表すと、

$$\Psi^{(1)} = \frac{\{\alpha(1-h)(1-P_A^{(0)}) + P_A^{(0)}\} \Psi}{\psi + P_A^{(0)}(1-\psi)} \quad (13)$$

よって

$$BW^{(1)} = 1 - (1 - \Psi^{(1)})^P + \alpha h \Psi^{(1)} P \quad (14)$$

以下同様に計算して、

$$BW_{local} = \lim_{k \rightarrow \infty} BW^{(k)} \quad (15)$$

を得る。

## 4.2 CGのモデル

一般にM個のノードがリング結合されている場合を考える。ノードのサイクルごとのリクエスト率を $\Psi$ 、ノードとリングのサイクルの比を $\beta$ とする。まず、リングのサイクルに対してi個のリクエストが発行される確率 $q(i)$ は

$$q(i) = \binom{M}{i} \left(\frac{\Psi}{\beta}\right)^i \left(1 - \frac{\Psi}{\beta}\right)^{M-i} \quad (16)$$

リングのサイクルにi個のリクエストがなされたとき、それらがアクセプトされる数を $E(i)$ とする。i個のリクエストのリクエスト先の総数は $M^i$ 。今、ある特定のノードはリクエストを受けないとすると、リクエストの総数は $(M-1)^i$ となるので、特定のノードが少なくとも1つリクエストを受ける確率は

$$1 - \left(\frac{M-1}{M}\right)^i \quad (17)$$

よって、M個のノードに対してi個のリクエストがなされたときにアクセプトされる確率は

$$E(i) = M - M \left(\frac{M-1}{M}\right)^i \quad (18)$$

仮定4で述べたように、リング・ネットワークの遅延時間は無視できないため、リクエストがアクセプトされる時間の近似を行わなければCGのバンド幅が推定できない。リングがアイドル状態ならリクエストはMサイクルで実行されるとする。CG内の各ノードでのリクエストの最大キューイ

ング数をLとすると、CG内でのリクエストの平均待ちサイクル数は

$$W = \sum_{i=0}^{LM} i \binom{M+i}{i} \left(\frac{\Psi}{\beta}\right)^i \left(1 - \frac{\Psi}{\beta}\right)^M \quad (19)$$

よって、CGでのバンド幅は

$$\frac{\sum_{i=0}^M E(i) q(i)}{M + W_{CG}} \quad (20)$$

式を簡単にすると

$$\frac{M}{M+W} \left\{ 1 - \left(1 - \frac{\Psi}{M\beta}\right)^M \right\} \quad (21)$$

## 4.3 クロスバ・ネットワークの分析

$K \times K$ のクロスバ・ネットワークの1入力に、リクエスト率rの入力が与えられたとする。K個の出力のうち特定の1出力にこのリクエストが伝えられる確率は $r/K$ である。特定の1出力にK個のどの入力のリクエストも伝えられない確率は

$$\left(1 - \frac{r}{K}\right)^K \quad (22)$$

よって、ある特定の出力に対して少なくとも1つの入力が伝えられる確率は

$$1 - \left(1 - \frac{r}{K}\right)^K \quad (23)$$

## 4.4 全体のモデル

P個のPCH付きプロセッサと1つのローカル・メモリ、NIFをバス結合したものをPC、M個のPCをリング結合したものをCG、K個のCGをクロスバ結合したものをASURAの現在の最大構成とする<sup>2</sup>。ただし、NIFには分散共有メモリの一部であるグローバル・メモリとGCHが含まれている。このグローバル・メモリはプロセッサからのアクセス時間によって、PCGM、CGGM、RGMの3階層に分れる。

システム全体のバンド幅BWは、ローカル・メモリ及びグローバル・キャッシュに対するバンド幅を $BW_{local}$ 、グローバル・メモリに対するバンド幅を $BW_{global}$ とすると、

$$BW = M K BW_{local} + K BW_{global} \quad (24)$$

で求められる。

プロセッサのローカル・メモリへのリクエスト率を $\Psi$ 、PCHのヒット率を $h_{PCH}$ 、GCHのヒッ

<sup>2</sup>現在のところ $P \times M \times K = 1024$

ト率を  $h_{GCH}$ 、プロセッサとバスのサイクル・タイムの比を  $\alpha$  とする。また、リクエスト先が LM、PCGM、CGGM、RGM である確率をそれぞれ、 $S_{LM}$ 、 $S_{PCGM}$ 、 $S_{CGGM}$ 、 $S_{RGM}$  とする。この時、

$$S_{LM} + S_{PCGM} + S_{CGGM} + S_{RGM} = 1 \quad (25)$$

が成り立つ。プロセッサから LM、GCH、PCGM、CGGM、RGM にそれぞれ実際に<sup>3</sup>リクエストされた確率、 $\Psi_{LM}$ 、 $\Psi_{GCH}$ 、 $\Psi_{PCGM}$ 、 $\Psi_{CGGM}$ 、 $\Psi_{RGM}$  は以下のように表される。

$$\Psi_{LM} = (1 - h_{PCH})S_{LM} \alpha \Psi \quad (26)$$

$$\Psi_{GCH} = (1 - h_{PCH})h_{GCH} \times (S_{PCGM} + S_{CGGM} + S_{RGM}) \alpha \Psi \quad (27)$$

$$\Psi_{PCGM} = (1 - h_{PCH})(1 - h_{GCH})S_{PCGM} \alpha \Psi \quad (28)$$

$$\Psi_{CGGM} = (1 - h_{PCH})(1 - h_{GCH})S_{CGGM} \alpha \Psi \quad (29)$$

$$\Psi_{RGM} = (1 - h_{PCH})(1 - h_{GCH})S_{RGM} \alpha \Psi \quad (30)$$

よってローカル・メモリ及びグローバル・キャッシュに対するバスへのリクエスト率は

$$\psi = \Psi_{LM} + \Psi_{GCH} + \Psi_{PCGM} \quad (31)$$

式 (2)(25) を利用して

$$BW_{local}^{(0)} = \{S_{LM} + S_{PCGM} + h_{GCH}(S_{CGGM} + S_{RGM})\} \times \{1 - (1 - \psi)^P\} + \alpha h_{PCH} \Psi P \quad (32)$$

以下 4.1 と同様に計算を行ない、

$$P_A^{(0)} = \frac{BW_{local}^{(0)}}{\alpha \Psi P} \quad (33)$$

$$\Psi^{(1)} = \frac{\psi(1 - P_A^{(0)}) + \Psi P_A^{(0)}}{P_A^{(0)} + \psi(1 - P_A^{(0)})} \quad (34)$$

$$BW^{(1)} = 1 - (1 - \Psi^{(1)})^P + \alpha h_{PCH} \Psi^{(1)} P \quad (35)$$

を繰り返すことにより

$$BW_{local} = \lim_{k \rightarrow \infty} BW_{local}^{(k)} \quad (36)$$

を得る。

CG のリング・ネットワークに PC から入力される確率を  $P_{CG}$ 、CG からクロスバ・ネットワークに入力される確率を  $P_X$  とする。 $\Psi_{CGGM}$  と  $\Psi_{RGM}$  は発行された PC の外へのリクエストであるから、

<sup>3</sup> PCH にヒットしないという意味

リングとバスとのサイクル・タイムの比を  $\beta$  とすると、PC 内の少なくとも 1 つのプロセッサが自分の PC 外にリクエストを出す確率、すなわち  $P_{CG}$  は、

$$P_{CG} = 1 - \left(1 - \frac{\Psi_{CGGM} + \Psi_{RGM}}{\beta}\right)^P \quad (37)$$

これが CG 内のリクエストである確率は

$$P_{CG} \frac{S_{CGGM}}{S_{CGGM} + S_{RGM}} \quad (38)$$

CG 外のリクエストである確率は

$$P_{CG} \frac{S_{RGM}}{S_{CGGM} + S_{RGM}} \quad (39)$$

よって、CG 内の少なくとも 1 つの PC が CG 外にリクエストを出す確率は

$$P_X = 1 - \left(1 - P_{CG} \frac{S_{RGM}}{S_{CGGM} + S_{RGM}}\right)^M \quad (40)$$

$P_X$  はクロスバ・ネットワークの 1 入力に与えられる確率で、伝えられるべき出力は自分以外の  $K-1$  個の出力のどれかである。式 (23) を利用すると、クロスバ・ネットワークの特定の出力に少なくとも 1 つの CG からの入力伝わる確率は、

$$1 - \left(1 - \frac{P_X}{K-1}\right)^{K-1} \quad (41)$$

これは特定の CG がクロスバ・ネットワーク経由でリクエストを受ける確率に等しいので、結局リングのリクエスト率  $\Psi_{ring}$  は

$$\Psi_{ring} = P_{CG} \frac{S_{RGM}}{S_{CGGM} + S_{RGM}} + 1 - \left(1 - \frac{P_X}{K-1}\right)^{K-1} \quad (42)$$

以上のことより、GC 内に  $i$  個のリクエストが発行される確率  $q_{CG}(i)$ 、それらがアクセプトされる確率  $E_{CG}(i)$ 、平均待ち時間  $W_{CG}$  は式 (16)(18)(19) よりそれぞれ

$$q_{CG}(i) = \binom{M}{i} \Psi_{ring}^i (1 - \Psi_{ring})^{M-i} \quad (43)$$

$$E_{CG}(i) = M - M \left(\frac{M-1}{M}\right)^i \quad (44)$$

$$W_{CG} = \sum_{i=0}^{LM} i \binom{M+i}{i} (\Psi_{ring})^i (1 - \Psi_{ring})^M \quad (45)$$

であるから、グローバル・メモリに関するバンド幅は式 (20) より

$$\frac{M}{M + W_{CG}} \left\{1 - \left(1 - \frac{\Psi_{ring}}{M}\right)^M\right\} \quad (46)$$

よって、全体のバンド幅は式(24)(15)(46)より

$$BW = K M BW_{local} + K BW_{global} \quad (47)$$

によって求められる。

## 5 評価

これまでに述べてきたモデルを用いて、ASURAの性能評価を行なう。性能評価の対象となるパラメータは、 $h_{PCH}$ と $S_{LM}$ である。ただし、 $h_{GCH} = 0.8h_{PCH}$ 、 $S_{CGGM} = (1 - S_{LM})/4$ 、 $S_{PCGM} = S_{CGGM}/4$ 、 $S_{RGM} = 1 - S_{PCGM} - S_{CGGM} - S_{RGM}$ という共通の設定を行なっている。

$h_{PCH}$ を固定する場合(0.9に固定) $S_{LM}$ の値を、 $S_{LM}$ を固定する場合(0.6に固定) $h_{PCH}$ の値を、それぞれ0から1まで0.01間隔に変えて評価を行なった。また、固定したパラメータは、 $\alpha = 1.0$ 、 $\beta = 1.0$ 、 $L = 5$ である。

図3はASURA全体のバンド幅を示している。 $h_{PCH}$ を固定した場合は $S_{LM}$ の値にあまり関係なく高いバンド幅を示している。また $S_{LM}$ を固定した場合、ほぼ $h_{PCH}$ の値に従った性能向上を示している。

図4はPCにおけるPCH、LM、GCH、PCGMのバンド幅を示している。図3とほとんど同じような性能向上を示していることから明らかに、ここの性能がシステム全体の性能を決定していると言えよう。

図5はCGにおけるCGGM、RGMのバンド幅を示している。パラメータの値によって、かなり複雑な結果を示しているが、全体に占めるバンド幅の量が極めて少ないため、システム全体の性能にはあまり影響を与えないと言えよう。

## 6 結論

分散共有メモリ型マルチプロセッサ「ASURA」の基本性能を、その階層性に着目した確率モデルで近似し、データの分散度合、及びキャッシュのヒット率が全体の性能にどのように影響を与えるかを確認した。その結果、全体の性能は主にキャッシュのヒット率に影響を受け、高ヒット率では理論性能に近い性能が可能であることを確認した。

本稿ではキャッシュのコヒーレント動作に関する定量化がなされなかったが、本モデルを拡張してキャッシュのコヒーレント動作を組み込んだ性能評価のモデルを開発すること、及び全体の正確なシミュレーションを行なうことが今後の研究課題である。

Overall Bandwidth

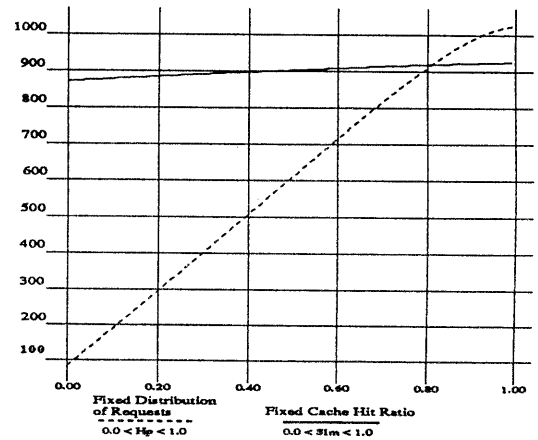


図3: ASURA全体の性能

PCH, LM, GCH Bandwidth

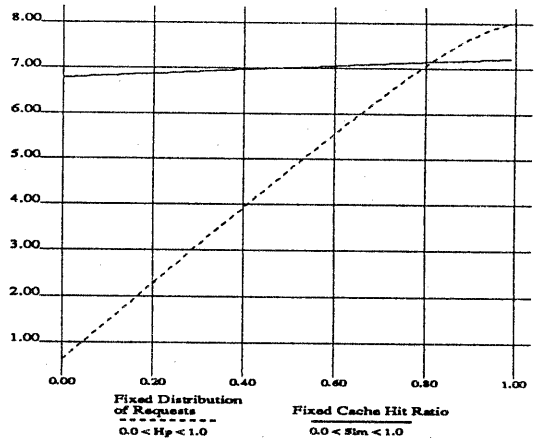


図4: キャッシュ、ローカル・メモリに関する性能  
謝辞

日頃ご討論頂く(株)クボタ コンピュータ事業推進室 山口宗之部長、及び同R&Dグループの諸氏に感謝致します。

## 参考文献

- [1] Santosh G. Abraham and Edward S. Davidson. A communication model for optimizing hierarchical multiprocessor systems. In *Proceedings of the 1986 International Conference on Parallel Processing*, pages 467-474, 1986.
- [2] Dharma P. Agrawal and Imadeldin O. Mahgoub. Performance analysis of cluster-based supersystems. In *Proceedings of the International Symposium on Computer Architecture*, pages 593-602, 1985.

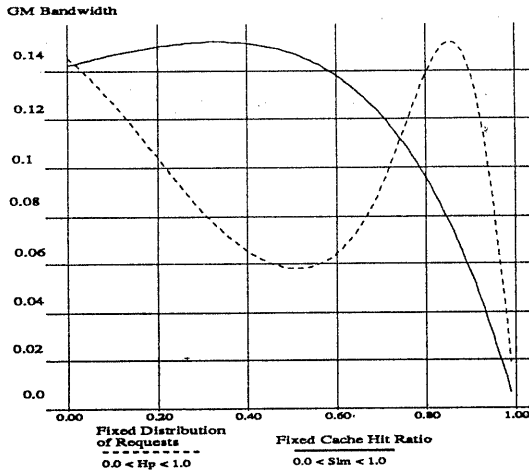


図 5: グローバル・メモリに関する性能

- [3] Dharma P. Agrawal and Imadeldin O. Mahgoub. Analysis of a class of cluster-based multiprocessor systems. *Information Science International Journal*, 43:85-105, 1987.
- [4] James Archibald and Jean-Loup Baer. Cache coherence protocols: Evaluation using a multiprocessor simulation model. *ACM Transactions on Computer Systems*, 4(4):273-298, 1986.
- [5] D. P. Bhandarkar. Analysis of memory interference in multiprocessors. *IEEE Transactions on Computers*, 24(9):897-908, 1975.
- [6] Laxmi N. Bhuyan. A combinatorial analysis of multibus multiprocessors. In *Proceedings of the 1984 International Conference on Parallel Processing*, pages 225-227, 1984.
- [7] Laxmi N. Bhuyan. An analysis of processor-memory interconnection networks. *IEEE Transactions on Computers*, 34(3):279-283, 1985.
- [8] Daniel Lenoski et al. The directory-based cache coherence protocol for the dash multiprocessor. In *Proceedings of the International Symposium on Computer Architecture*, pages 148-159, 1990.
- [9] David R. Cheriton et al. Paradigm: A highly scalable shared-memory multicomputer architecture. *IEEE Computer*, pages 33-46, 1991.
- [10] Manish Gupta and Prithviraj Banerjee. Demonstration of automatic data partitioning techniques for parallelizing compilers on multicomputers. *IEEE Transactions on Parallel and Distributed Systems*, 3(2):179-193, 1992.
- [11] Kai Hwang and Faye A. Briggs. *COMPUTER ARCHITECTURE AND PARALLEL PROCESSING*. McGraw-Hill, 1985.
- [12] David J. Kuck, Edward S. Davidson, and Duncan H. Lawrie Ahmed H. Sameh. Parallel supercomputing today and the cedar approach. *電子情報通信学会論文誌*, J71-D(8):1361-1374, 1988.
- [13] Monica S. Lam, Edward E. Rothberg, and Michael E. Wolf. The cache performance and optimizations of blocked algorithms. In *Proceedings of the Fourth International Conference on ASPLOS*, pages 63-74, 1991.
- [14] Tomas Lang, Mateo Valero, and Ignacio Alegre. Bandwidth of crossbar and multiple-bus connections for multiprocessors. *IEEE Transactions on Computers*, 31(12):1227-1234, 1982.
- [15] Imadeldin O. Mahgoub and Dharma P. Agrawal. Impact of cluster network failure on the performance of cluster-based supersystems. In *Proceedings of the 1986 International Conference on Parallel Processing*, pages 743-749, 1986.
- [16] Imadeldin O. Mahgoub and A. K. Elmagarmid. Performance analysis of a generalized class of m-level hierarchical multiprocessor systems. *IEEE Transactions on Parallel and Distributed Systems*, 3(2):129-138, 1992.
- [17] Marco Ajmone Marsan, Gianfranco Balbo, Gianni Conte, and Francesco Gregoretti. Modeling bus contention and memory interference in a multiprocessor system. *IEEE Transactions on Computers*, 32(1):60-72, 1983.
- [18] Marco Ajmone Marsan and Mario Gerla. Markov models for multiple bus multiprocessor systems. *IEEE Transactions on Computers*, 31(3):239-248, 1982.
- [19] T. N. Mudge and H. B. Al-Sadoun. A semi-markov model for the performance of multiple-bus systems. *IEEE Transactions on Computers*, 34(10):934-942, 1985.
- [20] T. N. Mudge, J. P. Hayes, and D. C. Winsor. Multiple bus architectures. *IEEE Computer Magazine*, 20(6):42-48, 1987.
- [21] Steven L. Scott, James R. Goodman, and Mary K. Vernon. Performance of the sci ring. In *Proceedings of the International Symposium on Computer Architecture*, pages 403-414, 1992.
- [22] Weijia Shang and Jose A. B. Fortes. Independent partitioning of algorithms with uniform dependencies. *IEEE Transactions on Computers*, 41(2):190-206, 1992.
- [23] Andrew S. Tanenbaum. *COMPUTER NETWORKS*. Prentice Hall, 1981.
- [24] Terry A. Welch. Memory hierarchy configuration analysis. *IEEE Transactions on Computers*, 27(5):408-413, 1978.
- [25] 齋藤 秀樹, 森 眞一郎, 城 和貴, David Fraser, 田中高士, and 富田 眞治. イベント対応型キャッシュ・コヒーレンス制御方式とそのバリア同期への応用. Technical Report 92-ARC-95-2, 情報処理学会計算機アーキテクチャ研究会, 1992.
- [26] 日本アライアントコンピュータ. *The CAM-PUS/800 Supercomputer*. 技術資料, 1991.
- [27] 森 眞一郎, 齋藤 秀樹, 五島 正裕, 富田 眞治, 田中高士, David Fraser, 城 和貴, and 新田 博之. 分散共有メモリ型マルチプロセッサ「阿修羅」の概要. Technical Report 92-ARC-94-6, 情報処理学会, 1992.